Communication-Efficient Decentralized Learning

Yuejie Chi

Carnegie Mellon University

EdgeComm Workshop, 2020

Acknowledgements



Communication-Efficient Distributed Optimization in Networks with Gradient Tracking and Variance Reduction, JMLR, 2020.

Distributed empirical risk minimization

Distributed/Federated learning: due to privacy and scalability, data are distributed at multiple locations / workers / agents.

Let $\mathcal{M} = \cup_i \mathcal{M}_i$ be a data partition with equal splitting:

$$f(\boldsymbol{x}) := \frac{1}{n} \sum_{i=1}^{n} f_i(\boldsymbol{x}), \quad \text{where} \quad f_i(\boldsymbol{x}) := \frac{1}{(N/n)} \sum_{\boldsymbol{z} \in \mathcal{M}_i} \ell(\boldsymbol{x}; \boldsymbol{z}).$$



$$\begin{array}{ll} {\sf minimize}_{\boldsymbol{x}} & f(\boldsymbol{x}) := \frac{1}{n}\sum_{i=1}^n f_i(\boldsymbol{x}) \\ & \downarrow \\ \\ {\sf minimize}_{\boldsymbol{x}_i} & \frac{1}{n}\sum_{i=1}^n f_i(\boldsymbol{x}_i) \quad {\sf subject to} \quad \boldsymbol{x}_i = \boldsymbol{x}_j \end{array}$$

$$\begin{array}{ll} {\sf minimize}_{\boldsymbol{x}} & f(\boldsymbol{x}) := \frac{1}{n} \sum_{i=1}^n f_i(\boldsymbol{x}) \\ & \Downarrow \\ \\ {\sf minimize}_{\boldsymbol{x}_i} & \frac{1}{n} \sum_{i=1}^n f_i(\boldsymbol{x}_i) \quad {\sf subject to} \quad \boldsymbol{x}_i = \boldsymbol{x}_j \end{array}$$

• Local computation: agents update local estimate;

 \Rightarrow need to be scalable!

m

$$\begin{array}{ll} {\rm minimize}_{\boldsymbol{x}} & f(\boldsymbol{x}) := \frac{1}{n}\sum_{i=1}^n f_i(\boldsymbol{x}) \\ & \Downarrow \\ \\ {\rm minimize}_{\boldsymbol{x}_i} & \frac{1}{n}\sum_{i=1}^n f_i(\boldsymbol{x}_i) \quad {\rm subject to} \quad \boldsymbol{x}_i = \boldsymbol{x}_j \end{array}$$

• Local computation: agents update local estimate;

 \Rightarrow need to be scalable!

• Global communications: agents exchange for consensus;

 \Rightarrow need to be communication-efficient!

$$\begin{array}{ll} {\rm minimize}_{\boldsymbol{x}} & f(\boldsymbol{x}) := \frac{1}{n}\sum_{i=1}^n f_i(\boldsymbol{x}) \\ & \Downarrow \\ \\ {\rm minimize}_{\boldsymbol{x}_i} & \frac{1}{n}\sum_{i=1}^n f_i(\boldsymbol{x}_i) \quad {\rm subject to} \quad \boldsymbol{x}_i = \boldsymbol{x}_j \end{array}$$

• Local computation: agents update local estimate;

 \Rightarrow need to be scalable!

• Global communications: agents exchange for consensus;

 \Rightarrow need to be communication-efficient!

Guiding principle: more local computation leads to less communication

Two distributed schemes



Master/slave model

PS coordinates <u>global</u> information sharing

Two distributed schemes





Master/slave model

PS coordinates <u>global</u> information sharing

Network model

agents share <u>local</u> information over a graph topology











Distributed Approximate NEwton (DANE) (Shamir et. al., 2014):

$$oldsymbol{x}_i^t = \operatorname*{argmin}_{oldsymbol{x}} f_i(oldsymbol{x}) - ig\langle
abla f_i(oldsymbol{x}^{t-1}) -
abla f(oldsymbol{x}^{t-1}), oldsymbol{x} ig
angle + rac{\mu}{2} \|oldsymbol{x} - oldsymbol{x}^{t-1}\|_2^2$$

Quasi-Newton method and less sensitive to ill-conditioning.



Distributed Stochastic Variance-Reduced Gradients (Cen et. al., 2020):

$$\boldsymbol{x}_{i}^{t,s} \Leftarrow \boldsymbol{x}_{i}^{t,s-1} - \eta \underbrace{\boldsymbol{v}_{i}^{t,s-1}}_{s}, \quad s = 1, 2, \dots,$$

• Better local computation efficiency.

Naive extension to the network setting



- **Communicate:** agent transmits $\{x_i^t, \nabla f_i(x_i^t)\}$;
- Compute:

$$\boldsymbol{x}_i^t \leftarrow \texttt{LocalUpdate}(f_i, \underbrace{\mathsf{Avg}\{\nabla f_j(\boldsymbol{x}_j^t)\}_{j \in \mathcal{N}_i}}_{\texttt{surrogate of } \nabla f(\boldsymbol{x}^t)}, \underbrace{\mathsf{Avg}\{\boldsymbol{x}_j^t\}_{j \in \mathcal{N}_i}}_{\texttt{surrogate of } \boldsymbol{x}^t})$$

Naive extension to the network setting



- Communicate: agent transmits $\{x_i^t, \nabla f_i(x_i^t)\}$;
- Compute:

$$\boldsymbol{x}_{i}^{t} \leftarrow \texttt{LocalUpdate}(f_{i}, \underbrace{\mathsf{Avg}\{\nabla f_{j}(\boldsymbol{x}_{j}^{t})\}_{j \in \mathcal{N}_{i}}}_{\texttt{surrogate of } \nabla f(\boldsymbol{x}^{t})}, \underbrace{\mathsf{Avg}\{\boldsymbol{x}_{j}^{t}\}_{j \in \mathcal{N}_{i}}}_{\texttt{surrogate of } \boldsymbol{x}^{t}}, \underbrace{\mathsf{Avg}\{\boldsymbol{x}_{j}^{t}$$

Naive extension to the network setting



- Communicate: agent transmits $\{x_i^t, \nabla f_i(x_i^t)\}$;
- Compute:

$$\boldsymbol{x}_i^t \leftarrow \texttt{LocalUpdate}(f_i, \underbrace{\texttt{Avg}\{\nabla f_j(\boldsymbol{x}_j^t)\}_{j \in \mathcal{N}_i}}_{\texttt{surrogate of } \nabla f(\boldsymbol{x}^t)}, \underbrace{\texttt{Avg}\{\boldsymbol{x}_j^t\}_{j \in \mathcal{N}_i}}_{\texttt{surrogate of } \boldsymbol{x}^t})$$

Consensus needs to be designed carefully in the network setting!

Average dynamic consensus

Assume that each agent generates some time-varying quantity r_i^t .

How to track its the dynamic average $\frac{1}{n}\sum_{j=1}^{n}r_{j}^{t} = \frac{1}{n}\mathbf{1}_{n}^{\top}\mathbf{r}^{t}$ in each of the agents, where $\mathbf{r}^{t} = [r_{1}^{t}, \cdots, r_{n}^{t}]^{\top}$?

Average dynamic consensus

Assume that each agent generates some time-varying quantity r_j^t .

How to track its the dynamic average $\frac{1}{n}\sum_{j=1}^{n}r_{j}^{t} = \frac{1}{n}\mathbf{1}_{n}^{\top}\mathbf{r}^{t}$ in each of the agents, where $\mathbf{r}^{t} = [r_{1}^{t}, \cdots, r_{n}^{t}]^{\top}$?

• Dynamic average consensus (Zhu and Martinez, 2010):

$$q^t = \underbrace{Wq^{t-1}}_{\mathsf{mixing}} + \underbrace{r^t - r^{t-1}}_{\mathsf{correction}},$$

where $\boldsymbol{q}^t = [q_1^t, \cdots, q_n^t]^\top$ and \boldsymbol{W} is the mixing matrix.

Average dynamic consensus

Assume that each agent generates some time-varying quantity r_j^t .

How to track its the dynamic average $\frac{1}{n}\sum_{j=1}^{n}r_{j}^{t} = \frac{1}{n}\mathbf{1}_{n}^{\top}\mathbf{r}^{t}$ in each of the agents, where $\mathbf{r}^{t} = [r_{1}^{t}, \cdots, r_{n}^{t}]^{\top}$?

• Dynamic average consensus (Zhu and Martinez, 2010):

$$q^t = \underbrace{Wq^{t-1}}_{\mathsf{mixing}} + \underbrace{r^t - r^{t-1}}_{\mathsf{correction}},$$

where $\boldsymbol{q}^t = [q_1^t, \cdots, q_n^t]^\top$ and \boldsymbol{W} is the mixing matrix.

• Key property: the average of $\{q_i^t\}$ dynamically tracks the average of $\{r_i^t\}$;

$$\mathbf{1}_n^{\top} \boldsymbol{q}^t = \mathbf{1}_n^{\top} \boldsymbol{r}^t,$$

M. Zhu and S. Martnez. "Discrete-time dynamic average consensus." Automatica 2010.

Gradient tracking

$$oldsymbol{x}_{i}^{t} \Leftarrow extsf{LocalUpdate}(f_{i},
and f_{i},
and f_{i},$$

• Parameter averaging:

$$\boldsymbol{y}_j^t = \sum_{k \in \mathcal{N}_j} w_{jk} \boldsymbol{x}_k^{t-1},$$

• Gradient tracking:

$$\boldsymbol{s}_j^t = \sum\nolimits_{k \in \mathcal{N}_j} w_{jk} \boldsymbol{s}_k^{t-1} + \underbrace{\nabla f_j(\boldsymbol{y}_j^t) - \nabla f_j(\boldsymbol{y}_j^{t-1})}_{\text{gradient tracking}}.$$

Gradient tracking

$$oldsymbol{x}_i^t \xleftarrow{} ext{LocalUpdate}(f_i,
and f(oldsymbol{x}^t),
oldsymbol{x}_j^t,
oldsymbol{x}_j^t$$

• Parameter averaging:

$$\boldsymbol{y}_j^t = \sum_{k \in \mathcal{N}_j} w_{jk} \boldsymbol{x}_k^{t-1},$$

• Gradient tracking:

$$\boldsymbol{s}_j^t = \sum\nolimits_{k \in \mathcal{N}_j} w_{jk} \boldsymbol{s}_k^{t-1} + \underbrace{\nabla f_j(\boldsymbol{y}_j^t) - \nabla f_j(\boldsymbol{y}_j^{t-1})}_{\text{gradient tracking}}.$$

We can now apply the same DANE and SVRG-type local updates!

Linear Regression: Well-Conditioned

$$f_i(\boldsymbol{x}) = \|\boldsymbol{y}_i - \boldsymbol{A}_i \boldsymbol{x}\|_2^2, \quad \boldsymbol{A}_i \in \mathbb{R}^{1000 \times 40}$$



Figure: The optimality gap w.r.t. iterations and gradients evaluation. The condition number $\kappa = 10$. ER graph (p = 0.3), 20 agents.

Linear Regression: Ill-Conditioned

$$f_i(\boldsymbol{x}) = \|\boldsymbol{y}_i - \boldsymbol{A}_i \boldsymbol{x}\|_2^2, \quad \boldsymbol{A}_i \in \mathbb{R}^{1000 \times 40}$$



Figure: The optimality gap w.r.t. iterations and gradients evaluations. The condition number $\kappa = 10^4$. ER graph (p = 0.3), 20 agents.

Extra Mixing

The mixing rate of the graph $\alpha_0 = 0.922$. A single round of mixing within each iteration cannot ensure the convergence of Network-SVRG.



Extra Mixing

The mixing rate of the graph $\alpha_0 = 0.922$. A single round of mixing within each iteration cannot ensure the convergence of Network-SVRG.



Extra Mixing

The mixing rate of the graph $\alpha_0 = 0.922$. A single round of mixing within each iteration cannot ensure the convergence of Network-SVRG.



Final remarks

- gradient tracking provides a way to extend master/slave algorithms (DANE and SVRG) to network settings;
- probing computational-communication trade-offs by employing different local updates and extra mixing;

Future work:

• convergence analysis in the nonconvex case.

Thank you!

Communication-Efficient Distributed Optimization in Networks with Gradient Tracking B. Li, S. Cen, Y. Chen, and Y. Chi, JMLR 2020.